

Tao Tao, Peter Cooper, Scott McGinnis, Wayne Matten, David Arndt, Greg Boratyn, and Tom Madden

Information Engineering Section, National Center for Biotechnology Information (NCBI), National Library of Medicine (NLM), National Institutes of Health (NIH)  
8600 Rockville Pike, Bethesda, Maryland 20814, USA

## Abstract

NCBI provides BLAST services through the BLAST homepage, where the public can search against the nucleotide database (NT), the non-redundant protein database (NR), a set of NCBI Reference Sequence Project (RefSeq) genomic sequences, and their annotated RNA or protein sequences. In this presentation, we will describe new features contained in the recently updated BLAST result pages that significantly enhance the page's usability. We describe the available RefSeq genomic databases and their relationship with genome records in NCBI's assembly database, which will enable more rational selection of databases and their taxonomic subset during a BLAST search. We will also cover alternative ways to access genomic sequences of interest, through assembly records on the web or FTP, for use in local standalone BLAST.

## NCBI Web BLAST Services

- Reside at its own domain ([blast.ncbi.nlm.nih.gov](http://blast.ncbi.nlm.nih.gov), 2006) with an reorganized homepage (2016, right)
- Offer easy access to core BLAST programs and a standard set of nucleotide and protein sequence collections as target databases
- Allow focused genomic data access through the BLAST Genomes search box
- Also provide access to specialized sequence analysis tools in "Specialized searches" section
- A well utilized set of services that serves over 50K users and does 350K interactive searches daily on average.

## Search Forms & Key Functions

### Databases and subset:

- Select the proper database from pulldown menu
- Specify a subset using Organism input & other preset filters (Entrez query no longer available for blastp)

### Search stringency settings:

- Upper limit of results via **Max target sequences**
- Significance setting via **Expect**
- Sensitivity via **Word size**
- Suppression of spurious hits through **Filter** and **Repeat** masking (nucleotide only)

Last one is often overlooked by our users!

## Updated Results Pages

- Places the **Search Summary** section at the top for convenient reference
- Brings out the **result filtering functions** from hiding for ready access
- Breaks results into **tabs** along the nature demarcation to reduce page scrolling
- Integrates Taxonomy report** as one of the tab sections [see right]
- Makes other **reformatting functions context-sensitive** to reduce confusion (more on this later)

## Filtering: an *E.coli* shiga toxin example

- Filter by **value ranges** (¥) to see desired subset of results satisfying the input criteria.
- Multiple ranges work in Boolean **AND** mode.

## Responsive Design

Functions appear where they matter & relevant

- Display format and CDS translation only in Alignment tab
- Sort HSPs only appears for subjects with multiple matches to a query
- Alignment length is only for pileup-like displays

- Filter by exclusion (activated by "exclude" checkbox, \$) to remove hits from the source organism that dominates the matches
- Filter by target organisms to see hits from interested sources. Multiple fields (added by "+" sign) operate in Boolean **OR** mode.

## Nucleotide Databases

The figure to the right summarizes ways data (both sequence and metadata) are organized and stored at NCBI. Understanding it will help we locate needed data more effectively.

These two refseq databases are in the standard set of BLAST databases accessible from search pages listed in the Web BLAST section of the BLAST homepage.

## Standard Nucleotide BLAST Databases

### General purpose & default

- NT**: Works fine for general purpose searches, but most eukaryotic genomes are NOT in it

### Transcripts

- Refseq\_rna**: Curated or annotated mRNA and rRNA from NCBI RefSeq

- TSA**: transcriptome assembly from submitters

- SRA**: RNA-seq

### Genomic Sequences

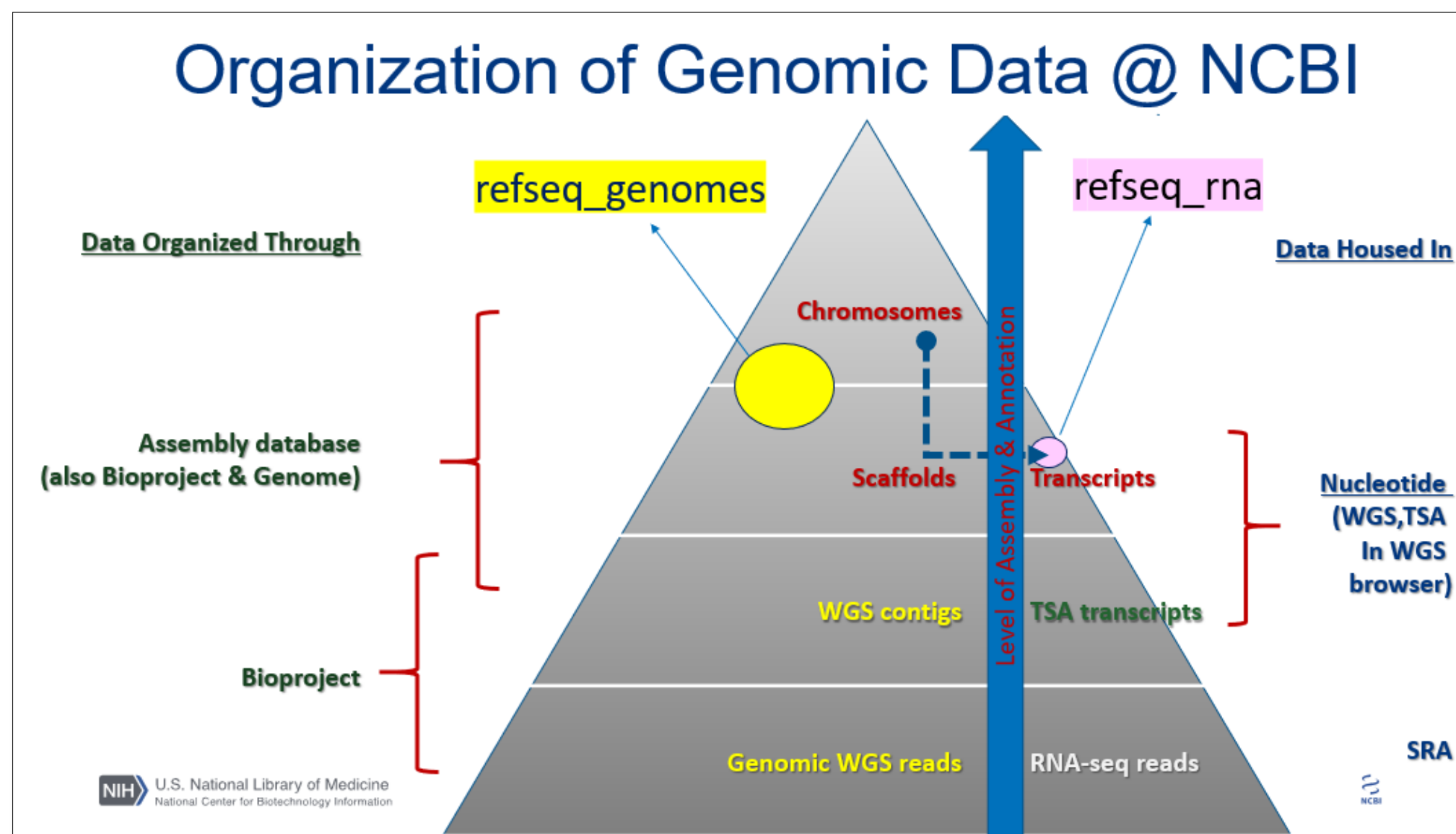
- Refseq\_genomic**: all genomic sequences (contigs, scaffolds, and full chromosomes), highly redundant

- Refseq\_genomes**: top-level genomic sequences, much less redundant

- Refseq\_representative\_genomes**: selective genomes from the above, mostly one set for a given organism (prokaryotes are exceptions)

- WGS**: whole genome shotgun contigs from submitters

- SRA**: Raw nextgen reads



## Genomic Data Available

The most informative way to get the summary and access the data is through the NCBI Eukaryotic Genome Annotation Pipeline's release list (partially shown above). The "B" link points to an organism-specific BLAST page, which provides access to the annotation release datasets, along with other related data for the organism.

For GenBank assemblies without refseq genome coverage, such as <https://go.usa.gov/xpASP>, their assembly records provide BLAST access to the original submitted assemblies. An easier way is to arrive at those organism specific genome BLAST page using the search box in the "BLAST Genomes" section of the BLAST homepage.

## Access Other less-Assembled Data

To access datasets from WGS, TSA, and SRA, use BLAST pages in the "Web BLAST" section, set WGS, TSA, or SRA as the database, then apply organism or project-based limit. Using "WGS Project" or "TSA Project" option and the letter initial of the project offers more reliable performance. For SRA, experiment or run access is the only option. Searches against experiment with many large runs (SRR) could be difficult to run/complete. Local data access through SRAtoolkit is recommended.

## Commandline Access to Genomic Data

This approach uses the **-remote** option of the standalone blast+ package. We need to know how to call the database of interest, and we can look up the database information through Entrez Programming Utilities or EntrezDirect:

\$ **esearch -db blastdbinfo -query 'ptaeda2.0' | esummary | xtract -pattern DocumentSummary -element Path genomic/3352/GCA\_000404065.3**

The whole string (value in the <Path> field) is the input argument for the -db option:

\$ **blastn -remote -db genomic/3352/GCA\_000404065.3 <...>**

## BLAST Access to Genomic Data Locally

Complete refseq\_genomic set (large redundant) is available for download from the BLAST ftp site. Use the update\_blastdb.pl to do this. Downloaded files are ready to use after extraction. Refseq\_genomes download is planned for the near future.

For organism-specific subset, we can refer to assembly record for ftp path to the \_genomic.fna.gz file for that assembly. We Need to format the sequence file into a blastable database using makeblastdb, before use. Batch downloading are described in the Assembly factsheet: [ftp.ncbi.nlm.nih.gov/pub/factsheets/Factsheet\\_Assembly.pdf](http://ftp.ncbi.nlm.nih.gov/pub/factsheets/Factsheet_Assembly.pdf)

## Local Data Access Using SRAtoolkit

Searching against WGS, TSA, and SRA datasets may be necessary for organisms without fully assembled genome, scaffold, or good transcript collections. NCBI's SRAtoolkit provides a convenient venue to search against these datasets under a local setup or in a rented cloud instance:

- prefetch the datasets of interests from NCBI

- blast search them locally using blastn\_vdb or tblastn\_vdb

A factsheets provides more details and dockerized packages are also available:

[ftp.ncbi.nlm.nih.gov/pub/factsheets/HowTo\\_Local\\_SRA\\_BLAST.pdf](http://ftp.ncbi.nlm.nih.gov/pub/factsheets/HowTo_Local_SRA_BLAST.pdf), [github.com/ncbi/docker](https://github.com/ncbi/docker)